

# Towards the mother-of-all-models: customised construction of the mark-recapture likelihood function

R. J. Barker & G. C. White

Barker, R. J. & White, G. C., 2004. Towards the mother-of-all-models: customised construction of the mark-recapture likelihood function. *Animal Biodiversity and Conservation*, 27.1: 177–185.

## Abstract

*Towards the mother-of-all-models: customised construction of the mark-recapture likelihood function.*— With a proliferation of mark-recapture models and studies collecting mark-recapture data, software and analysis methods are being continually revised. We consider the construction of the likelihood for a general model that incorporates all the features of the recently developed models: it is a multistate robust-design mark-recapture model that includes dead recoveries and resightings of marked animals and is parameterised in terms of state-specific recruitment, survival, movement, and capture probabilities, state-specific abundances, and state-specific recovery and resighting probabilities. The construction that we outline is based on a factorisation of the likelihood function with each factor corresponding to a different component of the data. Such a construction would allow the likelihood function for a mark-recapture analysis to be customized according to the components that are actually present in the dataset.

Key words: Capture-recapture, Mark-recapture, Likelihood.

## Resumen

*Hacia el origen de todos los modelos: una construcción personalizada de la función de verosimilitud en los estudios de marcaje-recaptura.*— Dada la proliferación de modelos de marcaje-recaptura y de los estudios que recopilan datos al respecto, los programas y los métodos de análisis se hallan sujetos a una continua revisión. En el presente estudio examinamos la construcción de la función de verosimilitud para un modelo general que incorpora todas las características de los modelos desarrollados recientemente. Se trata de un modelo multiestado robusto de marcaje-recaptura, que incluye las recuperaciones de individuos muertos y los reavistajes de animales marcados, parametrizándose con respecto al reclutamiento específico a un estado, la supervivencia, el movimiento y las probabilidades de captura, las abundancias específicas de un estado, las recuperaciones específicas de un estado y las probabilidades de reavistaje. La construcción que presentamos se basa en una factorización de la función de verosimilitud, de modo que cada factor corresponde a un componente distinto de los datos. Dicha construcción permitiría personalizar la función de verosimilitud en los análisis de marcaje-recaptura de acuerdo con los componentes que están realmente presentes en el conjunto de datos.

Palabras clave: Captura-recaptura, Marcaje-recaptura, Probabilidad.

Richard J. Barker, Dept. of Mathematics and Statistics, Univ. of Otago, P. O. Box 56, Dunedin, New Zealand.— Gary C. White, Dept. of Fishery and Wildlife Biology, Colorado State Univ., Fort Collins, Colorado, 80523 U.S.A.

## Introduction

The ability of biologists to collect varied and interesting data on the re-encounter of marked animals means that existing methods of analysis are often inadequate. This has been a major driver in the development of new mark-recapture models over the past 10–20 years and emphasises the importance of the relationship between data collection, models, and software development. As software has become increasingly available and mark-recapture models increasingly understood, more diverse mark-recapture data have become available. This in turn has spurred the development of new methods.

Our ability to learn from studies has always been limited by computing ability and the availability of suitable software. Early mark-recapture software such as program BROWNIE (Brownie et al., 1985) or JOLLY (Pollock et al., 1990) was relatively easy to use but offered relatively little choice in models. Three significant advances in mark-recapture software were: (1) program SURVIV (White, 1983) which allowed customized development of general multinomial models; (2) program SURGE (Lebreton & Clobert, 1986) which allowed use of a generalized linear model framework for open population mark-recapture models; and (3) program MARK (White & Burnham, 1999) which offered a "user-friendly" windows-based front-end that allows a wide choice of models and the flexibility of the generalized linear modelling approach. Program SURGE has also been generalized to MSURGE (Choquet et al., 2004) which allows multi-state models to be fitted with the model specified through a windows-type front-end.

In the past 10 years in particular there has been a proliferation of mark-recapture models. As these models have developed they have been incorporated into program MARK. This has culminated in a rapidly-lengthening model list of model choices confronting a user beginning a new analysis in MARK. At the 1997 EURING meeting the new analysis menu in MARK offered a choice between 10 models. Currently, there are 47 models for the user to choose from.

The models available in MARK can be categorized as open-population models that are conditional on release including the Cormack–Jolly–Seber (CJS) model (Cormack, 1964; Jolly, 1965; Seber, 1965), band recovery models (Brownie et al., 1985), multistate models (Brownie et al., 1993; Schwarz et al., 1993), joint recapture, dead recovery, and resighting models (Barker, 1997; Burnham, 1993), closed population models (Otis et al., 1978; Norris & Pollock, 1996; Pledger, 2000), robust design models (Pollock, 1981; Kendall et al., 1997) including robust-design generalizations of multistate models (Schwarz & Stobo, 1997; Kendall & Bjorkland, 2001) and of joint recapture, resighting and recovery models (Lindberg et al., 2001), and open population models that consider abundance or population growth (Pradel, 1996; Schwarz & Arnason, 1996). The correct choice of model depends on the nature of the

data and the parameters of interest. Good knowledge of the mark-recapture literature is required for the user to correctly select the model.

Many of the models in MARK represent small variations on a basic model type. For example the joint models of Burnham and Barker are closely related: Burnham's model extends the CJS model to include data from dead recoveries of animals and Barker's model extends Burnham's to include resightings of animals between sampling occasions. The Pradel "survival and seniority", "survival and lambda", and "survival and recruitment" models represent three different parameterisations of the same model. The "Recoveries only" and "Brownie at al. recoveries" models in MARK also represent different parameterisations of the same model. A strength of MARK is the ability to add constraints and covariate effects to the model easily in a generalized linear modelling framework.

Although the generalized linear model approach is very flexible, multiple parameterisations of the same model are included in MARK because constraints that are easily included under one parameterization may be difficult to include under another. For example, an analysis of open population data using Pradel's (1996) model in which recruitment rate  $f$  is equal for all periods, or follows a linear trend on the log-scale, is easily implemented if the model is parameterized in terms of  $\phi$  and  $f$ . Including this same constraint when the model is parameterized in terms of  $\lambda = \phi + f$  is more difficult. If a logit-link is used for  $\phi$  and a log-link for  $\lambda$  in this situation to model the effect of covariates, say, the required constraints become nonlinear and cannot be coded using the design matrix. Another example where constraints that are linear under one parametrization and nonlinear under another is the closed-population heterogeneity model of Agresti (1994) and Tjur (1982). Using a log-linear specification in which the encounter history probabilities are expressed in terms of parameters representing log-odds of capture and log-odds ratios the required constraints are linear. If the model is instead parameterised in terms of capture ( $p$ ) and recapture ( $c$ ) probabilities, as is done in MARK, the required constraints are nonlinear.

As the choice of models and complexity of software grows it is worth reflecting on whether the software can be improved. Below are some requirements we consider necessary for good mark-recapture software:

1. A choice of models that allows full use of all relevant data should be available. The choice of model should be guided in an obvious way based on the data available and the parameters of interest. One challenge is incorporating data from several sources. For example, a mark-recapture study might generate live recaptures, dead recovery, and resighting data. Alternatively, data might be available from several simultaneous studies. Multiple sources of data are important because it is often expensive to mark large numbers of animals but follow-up information can be relatively easy to

obtain. This is especially true as marking technology improves. Using radio-telemetry it is now possible to obtain large amounts of data, often in continuous time, from a relatively small number of tags. Making full use of these data is important so that maximum value is obtained from the study.

2. There should be flexibility to include nonlinear constraints in a simple way. It should also be possible for the user to estimate functions of parameters without having to reparameterize the model.

3. The software needs to have effective methods for dealing with large numbers of nuisance parameters. Many of the parameters in mark-recapture models are needed to describe aspects of the sampling process that are not of biologically intrinsic interest. The focus of the biologists should be on modelling interesting demographic patterns rather than the sampling processes required to obtain useful estimates of demographic parameters.

4. The models that are available should lend themselves to hierarchical extensions. Hierarchical modelling is of increasing interest as methods for fitting these models develop. Firstly, hierarchical models allow modeling of the relationship between parameters and parameter covariates; this process allows relevant biological questions to be answered. Secondly, they provide a framework for dealing with large numbers of nuisance parameters. Of the two reasons, the first is perhaps the most important. To illustrate this point, suppose it were possible for biologists to visit a study site regularly, and instead of marking and recapturing animals, they could somehow record the correct values of critical parameters. The biologists would almost certainly want to treat these repeated observations as data and fit some sort of model to the parameters. Hierarchical mark-recapture models that allow parameters to be modelled as random variables (i.e., as "data") are poorly developed. Most mark-recapture models focus on modeling the sampling process and regard the parameters as fixed constants to be estimated. In many respects, this focuses on the wrong process; it is the demographic mechanism that generated the parameters that should be of primary focus.

Recently developed models such as the joint model of Burnham (1993) or the robust design model (Kendall et al., 1997) have a specific likelihood function that has been coded in program MARK. We can envisage, however, a general model that incorporates all the features of the recently developed models: it is a multi-state robust-design mark-recapture model that includes dead recoveries and resightings of marked animals and is parameterised in terms of state-specific recruitment, survival and capture probabilities, state-specific abundances, and state-specific recovery and resighting probabilities.

One approach to developing software for a general model would be to write code that programs the complete likelihood function for the model. While it could be used as a basis for virtually all mark-recapture analyses it would be far too general for most studies. Irrelevant components of the model

would need to be disabled by making appropriate parameter constraints. An alternative approach would be to compartmentalize the model in such a way that, through an interface, the user was able to construct a customized likelihood by making appropriate selections from a choice of modules. Future work would focus on developing new modules rather than completely rewriting the likelihood function to accommodate new sampling and modeling developments.

In the rest of this paper we outline how such a customized likelihood function might be constructed.

#### A general likelihood for capture-recapture models in discrete-time

An underlying feature of almost every mark-recapture (MR) model is the need to model the capture process. Data from a MR experiment can be considered as repeated categorical measures with missing values. These missing values arise because animals can avoid capture and because the probability that an observation is missing depends on the survival status of the animal, it is important that the capture process is modelled correctly. Two approaches to modelling MR data have been developed. The first is based on direct modelling of the encounter histories using capture and survival probabilities. These approaches are exemplified by the sequence of models  $M_0$  through  $M_{tbh}$  as coded in program CAPTURE (Otis et al., 1978) for closed populations and by models such as the CJS model and the band return models of Brownie et al. (1985). These models are parameterised directly in terms of capture and survival probabilities. Model selection tends to focus on selecting a simplified description of the sampling process, and comparing different demographic parameterisations. The capture probabilities are often regarded as nuisance parameters in these analyses.

More recently, loglinear models have been used for analysing mark-recapture data. Most of the work has been on closed population models, although there has been some work on modelling open population models. The log-linear approaches pioneered by Fienberg (1972), Cormack (1989) and Agresti (1994) provide a general framework for analysing mark-recapture data. In this approach the data can be thought of as contributing to a  $2^k$  contingency table with each sample generating a binary classification according to whether or not an animal is caught (0,1). The cell corresponding to the null history 00...0 is missing. The likelihood function is expressed as a linear function on the log (or multinomial logit) scale of main-effects and interactions. The most general (saturated) model contains  $k$  main-effects and all interactions up to and including the  $(k-1)$ -way interaction. A  $k$ -way interaction cannot be fitted due to the fact that the null cell 00...0 is unobservable.

The emphasis in the log-linear approach is in selecting between models with different interaction terms included in the model. For closed populations,

model  $M_0$  corresponds to a model in which the main effects are equal, and all interactions are zero. Model  $M_1$  arises when the main effects are different but all interactions are zero. For these populations the interactions are usually regarded as nuisance parameters; the emphasis is on adopting a parsimonious model to obtain an estimate of abundance with small error. For closed populations the log-linear approach is a powerful technique that allows a full range of models to be fitted including versions of  $M_b$ ,  $M_{tb}$ ,  $M_h$ ,  $M_{bh}$  and  $M_{tbb}$  in a likelihood framework.

The loglinear approach has been applied to open populations by Cormack (1994) and more recently Rivest & Daigle (2004). For these populations, certain interactions must be included in the model to allow for mortality and births. The research emphasis has largely been on identifying equivalent log-linear models for standard models such as the Jolly-Seber model (Cormack, 1994) or the robust design (Rivest & Daigle, unpublished). However, there has been recognition of the potential that log-linear models have for increasing the richness of model structure in a parsimonious manner. For large studies involving many years the potential number of models is very large. For example, a 20-year single-state mark-recapture model has the potential for a model that has up to 1,048,575 parameters with a very much larger number of reduced parameter models! Clearly we are only ever likely to explore a very small fraction of these potential models, most of which will be too complicated to be useful. This large number of potential models also makes clear the importance of identifying a reasonable set of possible models before model fitting and that an approach based on a philosophy of exploring all possible is fruitless. An advantage of the log-linear approach is that the fit of the model is often improved by the addition of a small number of interaction terms that allow some dependencies between samples but without the large number of parameters needed to specify full dependency.

Although log-linear models have only been explored for closed populations, and for straight-forward Jolly-Seber and robust-design type models, they can in principle be extended to allow for other types of data including dead-recoveries and live-resightings, and to multiple capture/recapture states. For multi-state models, each additional state contributes another level within the capture classification. That is, instead of a  $2^k$  cross-classification it leads to a  $S^k$  cross-classification. The inclusion of dead recoveries and live resighting also contribute additional classification states. We can envisage then a general model that has multi-state captures/recaptures, dead recoveries, and live resightings. Instead of a simple 0 (not caught) and 1 (caught) we can extend the MARK LDLD...LD data classification to 00,10,20,...,S0,11,21,...,S1,12,22,...,S2,01,02 where the "L" member of the LD pair indicates capture state 0,1,...,S and the "D" member of the LD pair for the interval beginning

with sample  $i$  takes the value 1 if the animal is found dead in  $[i, i+1)$ , 2 if resighted alive in  $[i, i+1)$ , and 0 otherwise. For this model, we have a  $k^{3(S+1)}$  cross-classification with the potential for  $k^{3(S+1)}-1$  identifiable parameters. For example, a 5-period model with 3 states could have up to 499 parameters increasing to 3,999 for a 10-period study. This result emphasises the need for constraints that reduce the number of parameters in order to provide relatively simple descriptions of such complicated data sets.

In the log-linear approach it is relatively easy to find a parametric summary of important features of the data using standard linear modeling and model selection techniques. By incorporating the appropriate interactions for describing between-sample dependencies virtually any mark-recapture model based on a multinomial likelihood in discrete time can be accommodated. The difficulty lies in interpreting the coefficients of the fitted model. For the biologist, the parameters need to be expressed in terms of natural parameters such as survival and capture probabilities. If the log-linear approach is to develop, research needs to focus on identifying constraints on interaction terms that correspond to a reasonable set of constraints on the capture/recovery/resighting processes such as age-dependence, temporary trap response, and simple heterogeneity (Agresti, 1994; Tjur, 1982).

An alternative to the log-linear approach is to identify a useful set of constrained models by directly modelling the capture and demographic processes. This has been the traditional approach used in mark-recapture modelling. Below, we sketch out how the very general model described above could be constructed by adding factors that correspond to different components of the data. Such a construction would allow the likelihood function for a mark-recapture analysis to be customized according to the components that are actually present in the dataset.

#### Factorization of the CJS and JS models

Principal among open population capture-recapture models is the Cormack-Jolly Seber (CJS) model, first developed by Cormack (1964) and later extended by Jolly (1965) and Seber (1965). Under the assumption that all animals have the same probability of survival between sample  $i$  and  $i+1$ , and the same probability of capture in sample  $i$ , data in the CJS model can be summarized by an encounter history, with the count  $X_w$  denoting the number of animals that share the history  $w$ . The probability model used for inference is based on the joint distribution of the set of encounter histories which we denote by  $\{X_w\}$ .

In the CJS model we condition on the first release of each animal; implicit in this is the idea that although the numbers of animals released in each sample are usually random variables, their distribution is uninformative about the parameters of interest: the survival and capture probabilities.



Following Link & Barker (in press) we summarize the number of animals first caught at  $i$  by  $u_i$ . Excluding the null history  $0 = 00...0$ , we can write the joint distribution of the encounter history counts conditional on the first captures as  $[\{X_{w \neq 0} | \{u_i\}\}]$  which factors as:

$$[\{X_{w \neq 0} | \{u_i\}\}] = [\{X_{w \neq 0} | \{R_i\}\}] \times [\{R_i | \{u_i\}\}]$$

The first term,  $[\{X_{w \neq 0} | \{R_i\}\}]$  is the distribution of the encounter history counts conditional on the releases and so represents the CJS model. The second term  $[\{R_i | \{u_i\}\}]$  represents the distribution of the releases conditional on the first captures. It is needed to complete the description of the encounter histories but is usually of little interest. If there are no losses on capture, then  $[\{R_i | \{u_i\}\}] = 1.0$  and the CJS model is  $[\{X_{w \neq 0} | \{u_i\}\}]$ . If there are losses on capture then the term  $[\{R_i | \{u_i\}\}]$  can be factored out. In our discussion below we refer to  $[\{X_{w \neq 0} | \{u_i\}\}]$  as the CJS model; it is understood that if losses on capture are present, they are factored out of the model.

Where interest is in all demographic processes contributing to population dynamics, the CJS model is inadequate as the only demographic parameters are the survival probabilities. Information in the data about birth and abundance processes is contained in the distribution of the  $\{u_i\}$ . For example, the Jolly–Seber model is represented by multiplying the distribution  $[\{X_{w \neq 0} | \{u_i\}\}]$  by  $[\{u_i | \{U_i\}\}]$  where  $U_i$  denotes the number of unmarked animals in the population at the time of sample  $i$  (Seber, 1982). Thus, the Jolly–Seber model can be constructed from three distinct factors, each representing a distinct aspect of the sampling process.

An alternative extension of the CJS model developed by Schwarz & Arnason (1996) building on ideas from Crosbie & Manly (1985) is to multiply  $[\{X_{w \neq 0} | \{u_i\}\}]$  by  $[u_i | N] \times [\{u_i | \{u_i\}\}]$  where  $[u_i | N]$  represents the distribution of the number of unmarked animals caught during the study conditional on  $N$ , number of animals that were ever available for capture during the study. The term  $[\{u_i | \{u_i\}\}]$  represents the distribution of captures of unmarked animals conditional on  $u_i$ ; that is, it models the first capture of an animal given that it was ever caught. If the survival probabilities are constrained to one, and if all animals that were ever available for capture were present in the study population at the start of the experiment, then

$$[u_i | N] \times [\{u_i | \{u_i\}\}] \times [\{X_{w \neq 0} | \{u_i\}\}]$$

describes the distribution of the encounter histories for a closed population study with population size  $N$ .

Although the term  $[\{X_{w \neq 0} | \{u_i\}\}]$  is discussed above with reference to the CJS model it can be any appropriate distribution describing a set of encounter histories conditional on the first captures. Of the models currently in MARK, the term  $[\{X_{w \neq 0} | \{u_i\}\}]$  is sufficiently general to describe the following models: Recaptures only, Known fates, Recoveries only, Both (Burnham), Both (Barker),

and Brownie et al. Recoveries. The Schwarz and Arnason formulation thus provides a convenient way of adding birth and population growth modeling to most of the open population models in MARK, something currently only available for "Popan" (Schwarz & Arnason, 1996) and "Pradel" (Pradel, 1996) models. With appropriate modification to the distribution  $[\{u_i | \{u_i\}\}]$  the model  $[\{X_{w \neq 0} | \{u_i\}\}]$  could be the following models currently in MARK: Robust design, Multi–Strata Recaptures only, Barker Robust Design, Multi–Strata – Live and Dead Enc. All of these models can be extended to model recruitment or growth simply by adding to the likelihood the term  $[\{u_i | u_i\}]$  describing the distribution of the first captures conditional on capture at least once during the study. A complete description of the data requires the term  $[u_i | N]$  where  $N$  is the number of animals that were ever available for capture during the experiment. This term is included in the MARK "Popan" model however Link & Barker (in press) show that for open populations, this term is approximately ancillary to the estimable parameters in the model; that is, it contains virtually no information about the estimable parameters. This provides justification for omitting  $[u_i | N]$ , the approach taken by Pradel (1996). For closed populations the term must be included in order to estimate  $N$ .

Currently, program MARK is structured so that a distinct likelihood function is written for each model. However, the ability to write the models in terms of components of the sampling and demographic process as described above for the Jolly–Seber model suggests that it may be possible to customize the likelihood function to recognize different types of data and question. Future developments would not need to re–derive existing components, but would instead improve existing components or add new ones. We envisage that a user would specify the type and structure of data and indicate whether modeling of births and abundances was of interest. The appropriate  $[\{X_{w \neq 0} | \{u_i\}\}]$  and  $[\{u_i | u_i\}]$  terms would then be constructed based on option choices.

#### Construction of $[\{X_{w \neq 0} | \{u_i\}\}]$

The core component of the mark–recapture model is the term  $[\{X_{w \neq 0} | \{u_i\}\}]$ . It is here that the key information about the capture and survival processes is obtained from the data. This component can also be factored and provides the key to adding information from dead recoveries and resightings. The distribution of  $[\{X_{w \neq 0} | \{u_i\}\}]$  is the product of multinomial factors with indices  $\{u_i\}$  and probabilities  $\{Pr(X_w | u_i)\}$ . An animal with encounter history  $w$  that was first caught at  $i$ , contributes to  $[\{X_{w \neq 0} | \{u_i\}\}]$  through

$$Pr(\text{history} = w | \text{first caught at } i) = Pr(X_w | u_i).$$

If  $l$  indexes the sampling occasion when the animal was last caught, we can factor  $Pr(X_w | u_i)$  in the CJS model into three parts:

$$Pr(X_w|u_i) = \\ Pr(\text{survival from sample } i \text{ to } l | \text{first caught at } l) \\ \times Pr(\text{recaptures and releases} | \text{first caught at } i \text{ and survival to } l) \\ \times Pr(\text{not caught again} | \text{last caught in sample } l).$$

These terms also occur in joint live–dead models (Burnham, 1993; Barker, 1997; Barker et al., in press) with slight modification. For these models and for some animals, we have the additional information that they have survived beyond the last capture period, at least until the last capture occasion that defined the start of the interval in which they were last resighted alive or were found dead. If we index this occasion by  $l$ , and the sampling occasion when the animal was first released by  $j$ , then for the joint models we have the factorization:

$$Pr(X_w|u_i) = \\ Pr(\text{survival from sample } i \text{ to } l | \text{first caught at } l) \\ \times Pr(\text{recaptures and releases} | \text{first caught at } i \text{ and survival to } l) \\ \times Pr(\text{resightings} | \text{first caught at } i \text{ and survival to } l, \text{ recapture and releases}) \\ \times Pr(\text{encounter history after } l | \text{alive at time of sample } l).$$

This represents the CJS model with an additional factor for the resightings but with  $Pr(\text{encounter history after } l | \text{alive at time of sample } l)$  replacing  $Pr(\text{not caught again} | \text{last caught in sample } l)$ . This term incorporates the contribution made by a dead recovery or a live resighting in  $[l, l + 1)$ .

The term  $Pr(\text{recaptures and releases} | \text{first caught at } i \text{ and survival to } l)$  is the same as the equivalent term in the CJS model if there is no or random temporary emigration, but differs if there is permanent emigration. The distinction between the Barker model and the Burnham model is that the latter does not include the resighting terms and the recaptures term has been modified in the Barker model to allow Markovian temporary emigration.

We can write  $Pr(X_w|u_i)$  for the Barker or Burnham models as the CJS model multiplied by

$$Pr(\text{resightings} | \text{first caught at } i \text{ and survival to } l, \text{ recapture and releases})$$

$$\frac{Pr(\text{encounter history after } l | \text{alive at time of sample } l)}{Pr(\text{not caught again} | \text{last caught in sample } l)}$$

so, at least in principle, we can construct the Barker or Burnham models from a core CJS model by adding these factors.

#### Closed population models

To obtain closed population models, we need to add the factors  $[u_i | N]$  and  $[u_i | u_i]$  to the  $[X_{w \neq 0} | \{u_i\}]$  core that forms the CJS model. Because the population is closed we also need the constraints that the birth rates are all zero, the

survival probabilities one, and that all animal that were ever in the study population were present at the start of the study. The CJS parameterization of  $[X_{w \neq 0} | \{u_i\}]$  is sufficiently general to accommodate behavior effects as these can be coded using the design matrix. With an interface that allows constraints to be made on main effects and interactions for a multinomial logit link function, the heterogeneity models of Agresti (1994) and Tjor (1982) could also be fitted.

#### Generalization to robust design

To generalize the CJS core, or the Burnham & Barker joint models, to robust design versions we need to add a factor that describes captures and recaptures of animals during the sequence of samples that together define primary sampling occasion  $i$ . In the robust design, each of the capture occasions indexed by  $i$  is subdivided into secondary occasions. If we let  $p_{ij}$  denote the probability of capture in secondary occasion  $j$  ( $j = 1, \dots, J$ ) of primary occasion  $i$  then the probability  $p_i = 1 - (1 - p_{i1})(1 - p_{i2}) \dots (1 - p_{iJ})$  corresponds to the capture probability in the CJS model or the joint models of Burnham and Barker. We can write:

$$[\text{Secondary and primary recaptures}] = \\ [\text{Secondary captures} | \text{Primary captures}] \times \\ \times [\text{Primary captures}]$$

where  $[\text{Primary captures}]$  is the joint distribution of the primary capture events (animal is caught at least once in primary period  $i$ ) and is governed by the primary capture probabilities  $p_i$ . To complete the generalization we need to add the factor  $[u_i | u_i]$  which introduces the birth and growth parameters into the model if these are of interest. Kendall et al. (1995) used this approach to develop the likelihood function for the robust design extension of the CJS model except that instead of  $[u_i | u_i]$  they add the factor  $[u_i | \{U_i\}]$  as in the Jolly–Seber model.

#### Generalization to Multi-state models

As for the single-state case, the full multi-state likelihood can be factored into three distinct parts:

$$[X_w^M | N^{(i)}] = [u^{(i)} | N^{(i)}] \times [\{u_i^{(s)}\} | u^{(i)}] \times [X_{w \neq 0} | \{u_i^{(s)}\}]$$

where  $X_w^M$  denotes the number of animals with the multi-state encounter history  $w$ ,  $N^{(i)}$  is the number of animals that were ever present in the population during the study,  $u^{(i)}$  is the number of unmarked animals that were caught during the study and  $u_i^{(s)}$  is the number of unmarked animals caught in state  $s$  in sample  $i$ . The term  $[X_{w \neq 0}^M | \{u_i^{(s)}\}]$  is the joint distribution of the counts of the multi-state encounter histories conditional on the numbers first caught at each time and state. As such it represents any multi-state model conditional on the first captures including the Arnason–Schwarz (Schwarz et al., 1993)

model, the memory model of Brownie et al. (1993) and any extensions that include dead recoveries or live resightings. As for the single-state model the term  $[\{u_i^{(s)}\} | u_i^{(c)}]$  contains the useful information about changes in abundance and recruitment and provides the means for adding these components to the multi-state model.

The multi-state model written as above is not in the form of an expression involving a CJS core with a multi-state factor(s) added to the model. As an alternative we can write:

$$[\{X_w^M\} | N^{(c)}] = [u_i^{(c)} | N^{(c)}] \times [\{u_i^{(c)}\} | u_i^{(c)}] \times [\{X_{w \neq 0}^{(c)}\} | \{u_i^{(c)}\}] \times [X_{w \neq 0}^M | \{X_{w \neq 0}^{(c)}\}]$$

Here,  $X_w$  represents the encounter history collapsed across the states and so the first three factors are in the same form as the extended CJS model discussed above. However, each of these factors represent distributions that are complicated functions of state-specific capture, movement and survival probabilities. Whether this is a productive avenue to explore depends on how tractable these expressions turn out to be. The advantage would be that the full model can be built up from a CJS-type core, however the functions at the heart of this core will need to accept multi-state arguments.

## Discussion

The development of general software that allows biologists to focus on biological questions rather than battle with technological limitations of data analysis is important. We have outlined one approach to software development that may be useful in this context. Regardless of the approach, key issues that need to be addressed in future developments include the need for a suitable mechanism for incorporating alternative model parameterisations, multistate complications, handling large numbers of nuisance parameters, and ability to carry out hierarchical modelling.

A mechanism for offering alternative model parameterisations is important because biologists are often interested in estimating functions of parameters. The invariance property of the maximum likelihood estimator (Cox & Hinkley, 1974) means that it is easy to find the MLE and its variance for a function of parameters. A second reason for reparameterizing is that some parameter scales are more natural for introducing certain kinds of constraints. The programming problem is finding a compromise between having a wide choice of parameter transformations and the ease of use of the software. For example, nonlinear constraints on parameters can be programmed using Lagrange multipliers, as in POPAN (Arnason & Schwarz, 2002). However, computer code for incorporating nonlinear constraints in this manner would require the parsing of complicated algebraic expression making the programming problem difficult.

The multistate model is an exciting model that seems to present special difficulties (Lebreton & Pradel, 2002). A particular problem is that the likelihood function can have multiple-maxima, particularly when constraints are introduced into the model. This problem is exacerbated by the very large number of parameters for even moderately-sized multistate problems making it difficult to adequately explore the likelihood function graphically. An advantage of the approach that we have outlined of partitioning the likelihood function into components is that the likelihood function could be optimised sequentially. This is the approach taken by Schwarz & Arnason (1996) and if applied to the multi-state model might make it easier to explore the problem parts of the likelihood function.

A crucial issue is the management of a large number of nuisance parameters. It is not difficult to construct a model for data that has several hundred parameters, many of which are of little interest to the researcher. This problem is exacerbated if unobservable and misclassified states are included in the model (Kendall, 2004). These nuisance parameters are needed however to correctly specify the model and to maximize the information that is extracted from the data about the demographic parameters. Because having many nuisance parameters causes reduced precision of parameter estimates it is usually desirable to decrease their number. The standard approach in mark-recapture modeling is to use model-selection to reduce the number of nuisance parameters. The main difficulty with this approach is that there are often a large number of plausible models with differing structure imposed on nuisance parameters. The need to carry out model-selection over these parameters is a distraction from what should be the core focus of model selection: exploring biologically interesting hypotheses by comparing a small set of models that have various restrictions on the demographic parameters. An alternative could be to fit a very general model with few restrictions on nuisance parameters and then consider biological hypotheses by restricting the demographic parameters. An intermediate approach is the use of hierarchical models in which nuisance parameters are modelled as random effects; currently software that allows this approach is only available for relatively simple models such as the CJS model.

An alternative to reducing the number of nuisance parameters by restricting them is to eliminate them entirely from the model. The elimination of nuisance parameters is a well-known and difficult problem (see Berger et al., 1999 for a review). This can be done by integrating the nuisance parameters out of the model; this is essentially a Bayesian approach and can be done with respect to a non-informative prior distribution for the nuisance parameters. Computationally this approach is prohibitive for anything but very simple mark-recapture models. Conditional likelihood can also be used to eliminate nuisance parameters. This approach is suitable when components of the suffi-

cient statistics have conditional distributions that depend only on the parameters of interest. The simplest example where conditional likelihood is used in mark-recapture analysis is the 2-sample closed-population model  $M_c$ . Assuming that the number of animals with the four possible capture histories 11, 10, 01, and 00 are multinomial with index  $N$  we can write:

$$\begin{aligned} [X_{11}, X_{10}, X_{01}, X_{00} | N] &= \\ &= \frac{\binom{n_1}{m_2} \binom{N-n_1}{n_2-m_2}}{\binom{N}{n_2}} \times \prod_{i=1}^2 p_i^{n_i} (1-p_i)^{N-n_i} \end{aligned}$$

where the sufficient statistics are  $n_i$ , the number of animals caught in sample  $i$  and  $m_2$ , the number of marked animals caught in sample 2. The hypergeometric term represents the distribution  $[m_2 | N]$  and does not depend on the nuisance parameters hence can be used as a conditional likelihood. Conditional likelihoods have also been used for modelling heterogeneity in closed populations (Agresti, 1994; Tjur, 1982) and in the memory model of Brownie et al. (1993). Conditional likelihoods remain an interesting subject for research but at present there is little scope for their use in mark-recapture modelling.

For a variety of reasons, including those outlined above, we believe that hierarchical models are going to become of increasing importance to biologists. Programs such as MARK, POPAN, SURGE and SURVIV, and the possible extensions that we have outlined above code the likelihood function for specific situations. Constructing likelihood functions for hierarchical models is prohibitive in most cases because of the multi-dimensional integrations required. An alternative approach that has recently become popular is vague prior Bayesian analysis based on MCMC, for example using program WinBUGS (Spiegelhalter, 2000). In this approach a prior distribution for parameters is specified, and this distribution itself may have a so-called hyperprior distribution. Inference is made by summarizing the posterior distribution which is approximated by Monte Carlo simulation. Simple models such as the CJS model are easily coded in WinBUGS, but it becomes more difficult for some of the more complicated models. An alternative to WinBUGS is to develop specific mark-recapture code. For example, a HyperMARK module could be developed that clips on a user-specified prior distribution to the distributions used for expressing the mark-recapture likelihood function. MCMC could then be used to sample from to sample from the posterior distribution. The key issues here are (1) writing the HyperMARK code in an efficient and user-friendly manner that makes specification of biologically relevant models easy and (2) the user would need to accept the logic of Bayesian inference. If we can accept a vague

prior as a statement of knowledge about parameters before the experiment then the Bayesian logic is impeccable, however the appropriate description of near-ignorance before an experiment is carried out is controversial.

## References

- Agresti, A., 1994. Simple capture-recapture models permitting unequal catchability and variable sampling effort. *Biometrics*, 50: 494–500.
- Arnason, A. N. & Schwarz, C. J., 2002. POPAN-6: exploring convergence and estimate properties with SIMULATE. *Journal of Applied Statistics*, 29: 649–668.
- Barker, R. J., 1997. Joint modelling of live-recapture, tag-resight, and tag-recovery data. *Biometrics*, 53: 666–677.
- Barker, R. J., Burnham, K. P. & White, G. C., in press. Encounter history modelling of joint mark-recapture, tag-resighting and tag-recovery data under temporary emigration. *Statistica Sinica*.
- Berger, J. O., Brunero, L. & Wolpert, R. L., 1999. Integrated Likelihood Methods for Eliminating Nuisance Parameters. *Statistical Science*, 14: 1–28.
- Brownie, C., Anderson, D. R., Burnham, K. P. & Robson, D. S., 1985. *Statistical Inference from Band-Recovery Data: A handbook*, 2<sup>nd</sup> edition. U.S. Fish and Wildlife Service, Resource Publication 156. Washington D.C., U.S. Dep. of the Interior.
- Brownie, C., Hines, J. E., Nichols, J. D., Pollock, K. H. & Hestbeck, J. B., 1993. Capture-recapture studies for multiple strata including non-Markovian transitions. *Biometrics*, 49: 1173–1187.
- Burnham, K. P., 1993. A theory for combined analysis of ring recovery and recapture data. In: *Marked Individuals in Bird Population Studies*: 199–213 (J.-D. Lebreton & P. North, Eds.). Birkhauser Verlag, Basel.
- Choquet R., Reboulet, A.-M., Pradel, R., Gimenez, O. & Lebreton, J. D., 2004. M-SURGE New software for multistate recapture models. *Animal Biodiversity and Conservation*, 27: 1.
- Cormack, R. M., 1964. Estimates of survival from the sighting of marked animals. *Biometrika*, 51: 429–438.
- 1989. Log-linear models for capture-recapture. *Biometrics*, 45: 395–413.
- 1994. Unification of mark-recapture analyses by loglinear modelling. In: *Statistics in Ecology and Environmental Monitoring* (D. J. Fletcher & B. F. J. Manly, Eds.). Otago Conference Series 2, Univ. of Otago Press, Dunedin.
- Cox, D. R. & Hinkley, D. V., 1974. *Theoretical Statistics*. Chapman and Hall, London, U.K.
- Crosbie, S. F. & Manly, B. F. J., 1985. Parsimonious modelling of capture-mark-recapture studies. *Biometrics*, 41: 385–398.
- Fienberg, S. E., 1972. The multiple-recapture census for closed populations and incomplete contingency tables. *Biometrika*, 59: 591–603.



- Jolly, G. M., 1965. Explicit estimates from capture–recapture data with both death and immigration stochastic model. *Biometrika*, 52: 225–247.
- Kendall, W. L., 2004. Coping with unobservable and misclassified states in capture–recapture studies. *Animal Biodiversity and Conservation*, 27.1: 97–107.
- Kendall, W. L. & Bjorkland, R., 2001. Using open robust design models to estimate temporary emigration from capture–recapture data. *Biometrics*, 57: 1113–1122.
- Kendall, W. L., Nichols, J. D. & Hines, J. E., 1997. Estimating temporary emigration and breeding proportions using capture–recapture data with Pollock's robust design. *Ecology*, 78: 563–578.
- Kendall, W. L., Pollock, K. H. & Brownie, C., 1995. A Likelihood-based approach to capture–recapture estimations of demographic parameters under the robust design. *Biometrics*, 51: 293–308.
- Lebreton, J.-D. & Clobert, J., 1986. *Users manual for program SURGE. Version 2.0*. C.E.F.E., C.N.R.S., Montpellier, France.
- Lebreton, J.-D. & Pradel, R., 2002. Multistate recapture models: modelling incomplete individual histories. *Journal of Applied Statistics*, 29: 353–369.
- Lindberg, M. S., Kendall, W. L., Hines, J. E. & Anderson, M. G., 2001. Combining band recovery data and Pollock's robust design to model temporary and permanent emigration. *Biometrics*, 57: 273–281.
- Link, W. A. & Barker, R. J. (in press) Hierarchical models for open population capture–recapture data. *Biometrics*.
- Norris, J. L. & Pollock, K. H., 1996. Non-parametric MLE under two closed capture–recapture models with heterogeneity. *Biometrics*, 52: 639–649.
- Otis, D. L., Burnham, K. P., White, G. C. & Anderson, D. R., 1978. Statistical inference from capture–recapture data in closed animal populations. *Wildlife Monographs*, 62: rates. *Biometrics*, 37: 521–529.
- Pollock, K. H., Nichols, J. D., Brownie, C. & Hines, J. E., 1990. Statistical inference for capture–recapture experiments. *Wildlife Monographs*, 107: 1–97.
- Pledger, S., 2000. Unified maximum likelihood estimates for closed capture–recapture models using mixtures. *Biometrics*, 56: 434–442.
- Pradel, R., 1996. Utilization of capture–mark–recapture for the study of recruitment and population growth rate. *Biometrics*, 52: 703–709.
- Rivest, L.-P. & Daigle, G., 2004. Loglinear models for the robust design in mark–recapture experiments. *Biometrics*, 60: 100–107.
- Schwarz, C. J. & Arnason, A. N., 1996. A general methodology for the analysis of capture–recapture experiments in open populations. *Biometrics*, 52: 860–873.
- Schwarz, C. J. Schweigert, J. F. & Arnason, A. N., 1993. Estimating migration rates using tag–recovery data. *Biometrics*, 49: 177–193.
- Schwarz, C. J. & Stobo, W. T., 1997. Estimating temporary migration using the robust design. *Biometrics*, 53: 178–194.
- Seber, G. A. F., 1965. A note on the multiple recapture census. *Biometrika*, 52: 249–259.
- 1982. *The Estimation of Animal Abundance and Related Parameter*. Charles Griffin and Co., London.
- Spiegelhalter, D., Thomas, A. & Best, N. G., 2000. *WinBUGS Version 1.3 user manual*. Medical Research Council Biostatistics Unit, Cambridge.
- Tjur, T., 1982. A connection between Rasch's Item Analysis Model and a Multiplicative Poisson Model. *Scandinavian Journal of Statistics*, 9: 23–30.
- White, G. C., 1983. Numerical estimation of survival rates from band–recovery and biotelemetry data. *Journal of Wildlife Management*, 47: 716–728.
- White, G. C. & P. Burnham, K. P., 1999. Program MARK: survival estimation from populations of marked animals. *Bird Study*, 46(supplement): 120–139.